

Data Mining Lab 2020: Variational Fair Autoencoder

Intermediate Results

Mohammad Abir, Fabian Beringer

VFAE Short Refresher

- encodes features into latent space \mathbf{z} using NNs
- naturally disentangles *sensitive variable* \mathbf{s} from \mathbf{x} by modeling them as independent
- also injects label information about \mathbf{y}
 - > keeps \mathbf{z} useful when predicting \mathbf{y} but removes \mathbf{s}



Dataset: Adult Income

Dataset

- Each instance **x** represents a person
- Labels **y**: whether the person has >50k annual income
- The protected attribute **s** is **gender**:
67% Male, 33% Female (protected class)
- 45,222 instances, 15 features each



Dataset

For pre-processing we:

- removed NaN's
- binarized each feature using one-hot encoding
- bucketized continuous features into **5** buckets each
- resulted in **117** binary features



Challenges

Paper Omitting Details

- pre-processing not explained in detail
 - we followed another paper that was mentioned as reference
- no info about total training times
 - we trained for 100 epochs at max
- penalty scaling parameter β unclear for some datasets:
“scaled according to validation set”
 - we found that $\beta=1.0$ i.e. not scaling at all worked best



Further Challenges

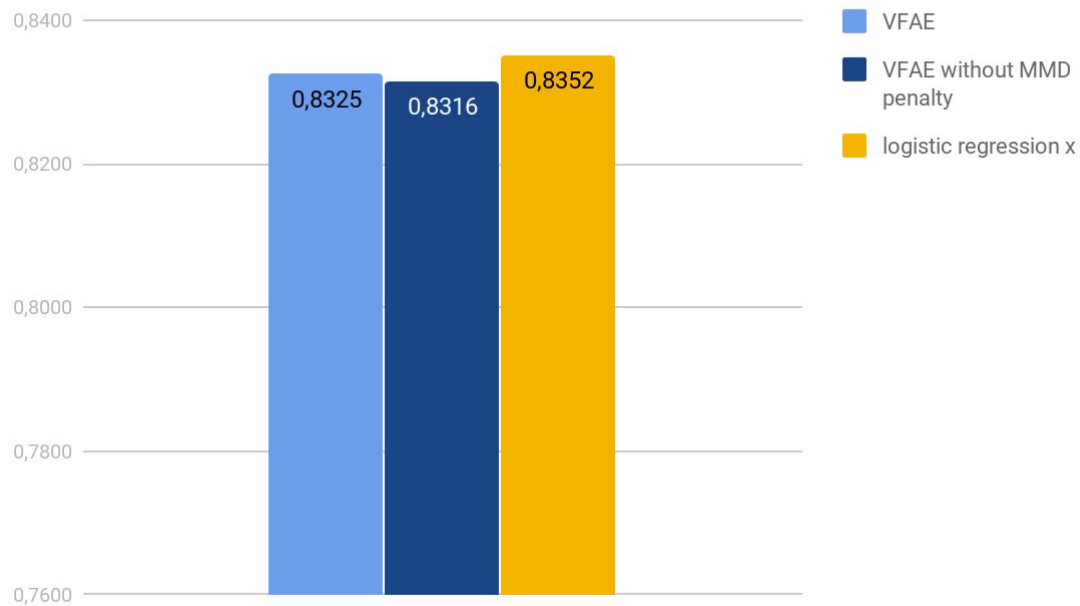
- rather complex objective function, harder to debug
- non standard metrics i.e. discrimination have some specialties when it comes to implementation details
- designing a training pipeline for different models which depend on each other can be a bit tricky



Comparing Results

Accuracy on y , trained on x vs. z

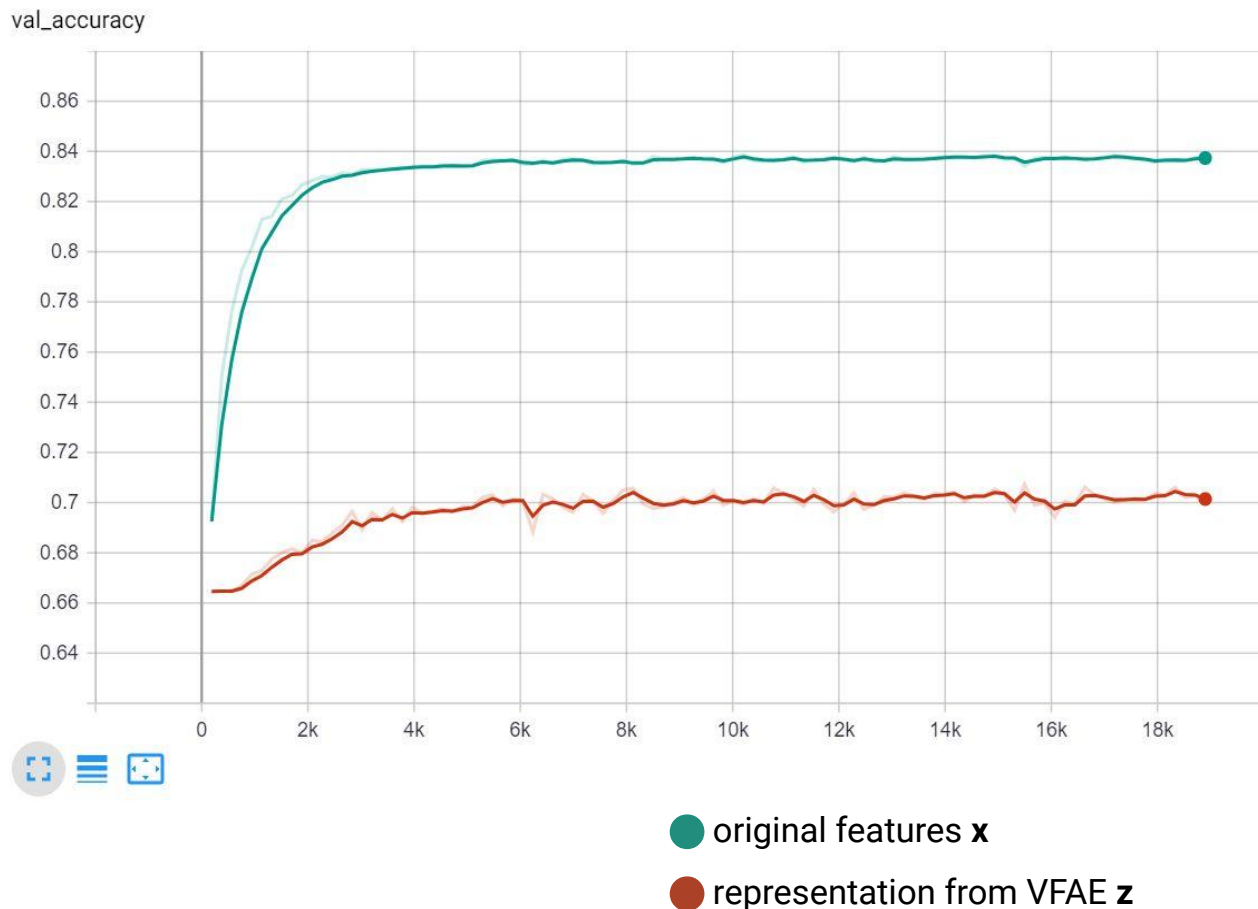
Accuracy Y (Income $\leq 50k$)



Accuracy on s

Validation accuracy during training when predicting gender.

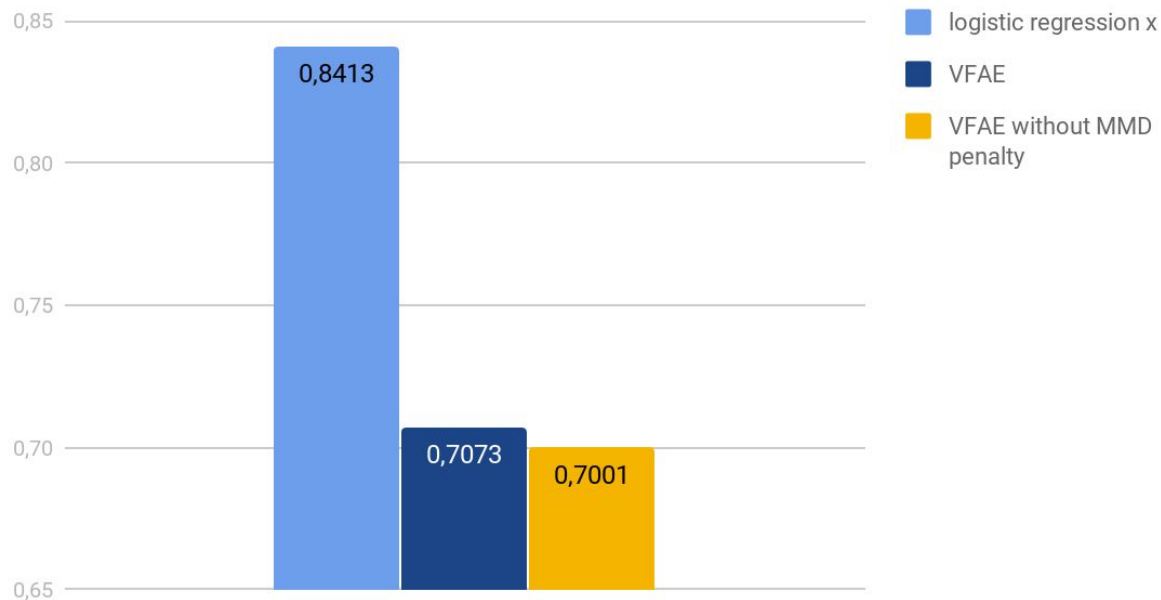
- slightly higher than mentioned in the paper
-> due to longer training time?



Accuracy on \mathbf{s} , trained on \mathbf{x} vs. \mathbf{z}

Accuracy Gender

test accuracy predicting gender (lower is better)



Discrimination Score

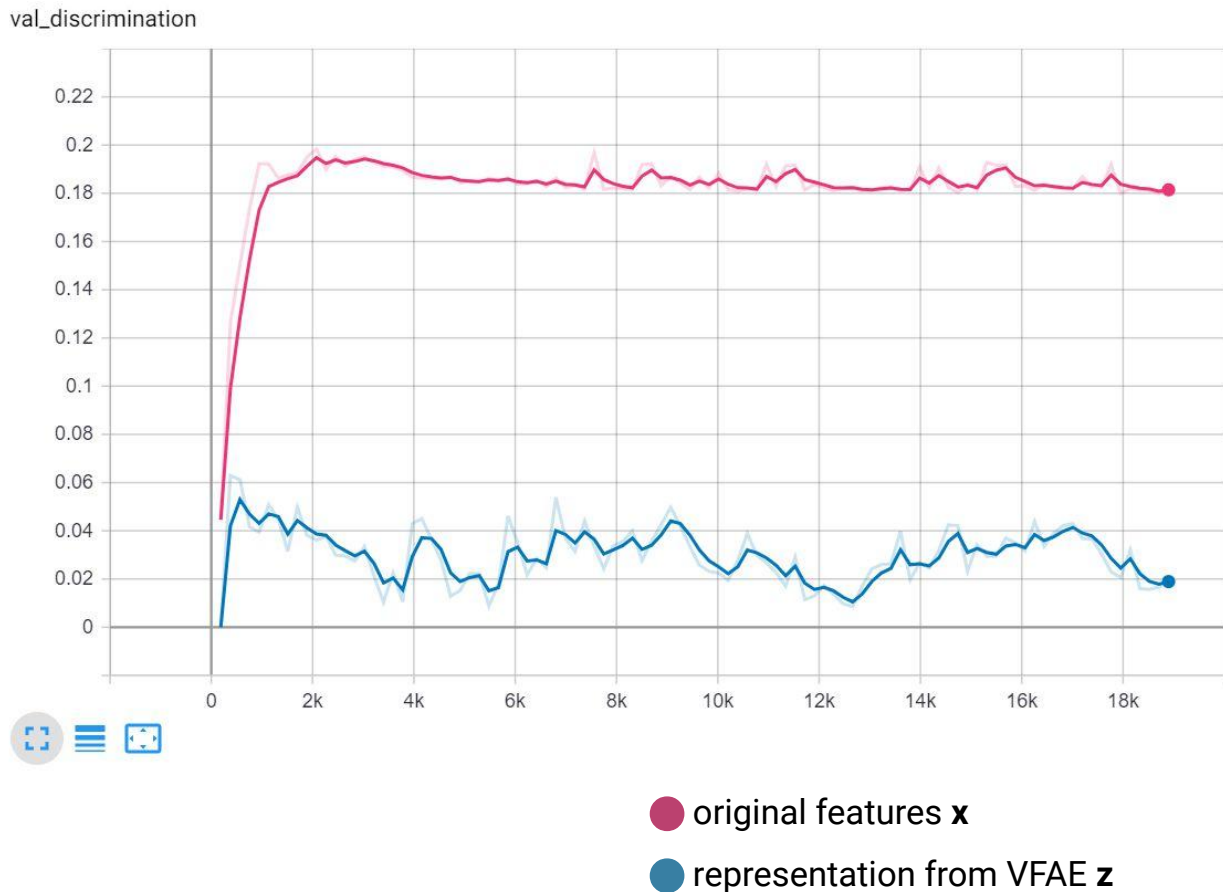
$$\text{Discrimination} = \left| \frac{\sum_{n=1}^N \mathbb{I}[y_n^{s=0}]}{N_{s=0}} - \frac{\sum_{n=1}^N \mathbb{I}[y_n^{s=1}]}{N_{s=1}} \right|$$

$$\text{Discrimination prob.} = \left| \frac{\sum_{n=1}^N p(y_n^{s=0})}{N_{s=0}} - \frac{\sum_{n=1}^N p(y_n^{s=1})}{N_{s=1}} \right|$$

Discrimination Score

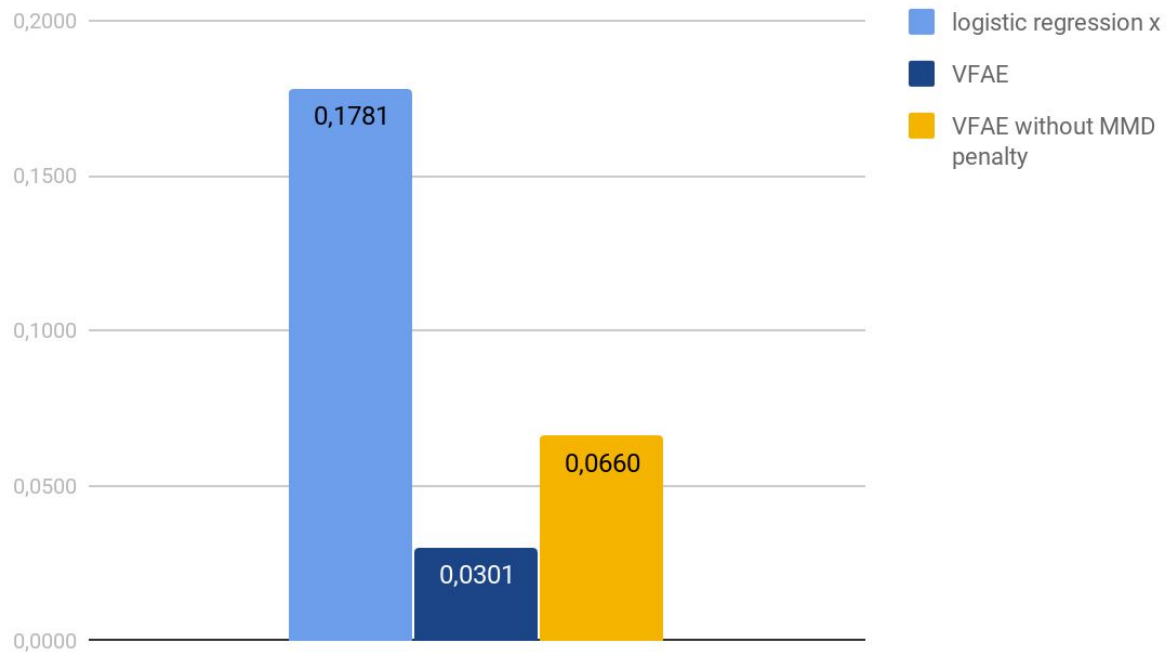
Discrimination on the validation set during training.

- discrimination on \mathbf{z} stays low rather consistently
→ gender is being factored out properly
- however \mathbf{z} is more noisy



Discrimination Gender

Discrimination Score



Effects of the MMD penalty

We found the use of the MMD penalty to

- result in higher accuracy on \mathbf{s} (lower is better)
- cause high variance in discrimination during training
- make it more likely to jump into good minima

-> room for improvement

$$\ell_{\text{MMD}}(\mathbf{Z}_{1\mathbf{s}=0}, \mathbf{Z}_{1\mathbf{s}=1}) = \left\| \mathbb{E}_{\tilde{p}(\mathbf{x}|\mathbf{s}=0)} [\mathbb{E}_{q(\mathbf{z}_1|\mathbf{x},\mathbf{s}=0)} [\psi(\mathbf{z}_1)]] - \mathbb{E}_{\tilde{p}(\mathbf{x}|\mathbf{s}=1)} [\mathbb{E}_{q(\mathbf{z}_1|\mathbf{x},\mathbf{s}=1)} [\psi(\mathbf{z}_1)]] \right\|^2$$



Predicting \mathbf{s} without MMD

- MMD penalty resulting in higher accuracy on \mathbf{s}
- could be due to:
 - different scaling
 - longer training time

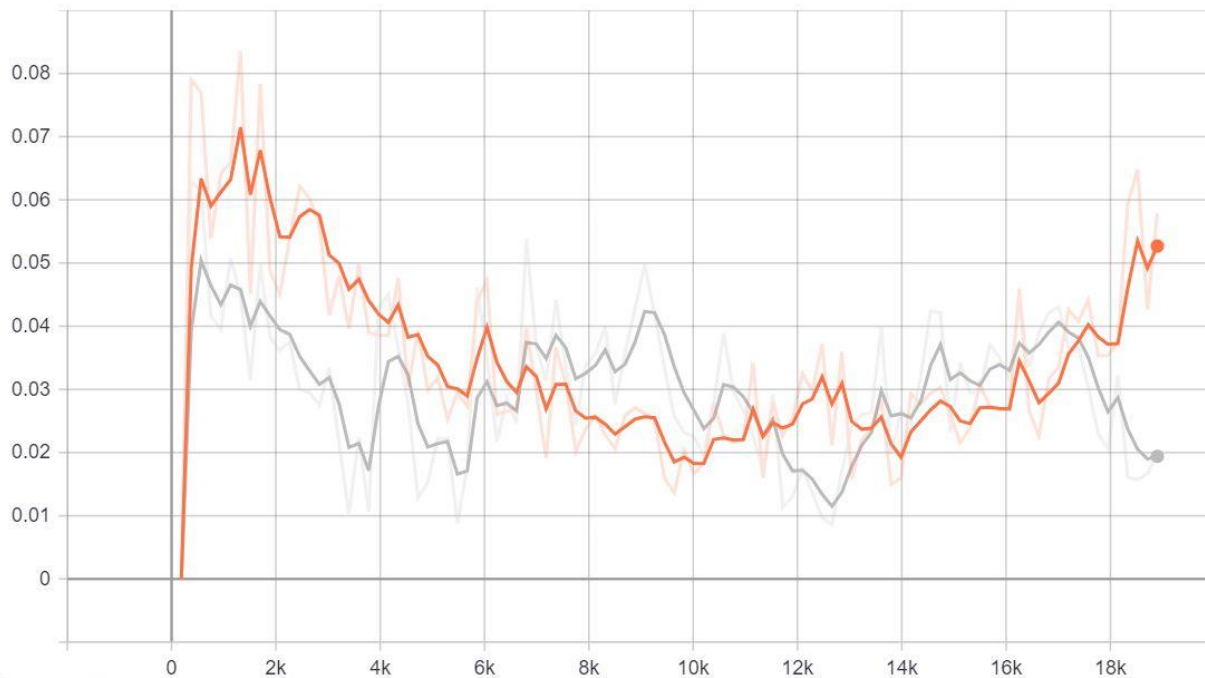


● with MMD penalty

● without MMD penalty

Discrimination with MMD more noisy

val_discrimination



Improvement Ideas

Improvement Ideas

- actual MMD instead of the fast approximation
- replace MMD with a different penalty
- try common ideas for AutoEncoders:
 - add input noise for robustness
 - use Contractive AutoEncoder penalty





Thanks for listening!